



## Review of recent Methodological Developments in group-randomized trials: Part 1 - Design

Mélanie Prague, Elizabeth Turner, Li Fan, Gallis John, Murray David

### ► To cite this version:

Mélanie Prague, Elizabeth Turner, Li Fan, Gallis John, Murray David. Review of recent Methodological Developments in group-randomized trials: Part 1 - Design. American Journal of Public Health, American Public Health Association, 2017. hal-01579073

**HAL Id: hal-01579073**

**<https://hal.inria.fr/hal-01579073>**

Submitted on 30 Aug 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **REVIEW OF RECENT METHODOLOGICAL DEVELOPMENTS IN**

2 **GROUP-RANDOMIZED TRIALS: PART 1 - DESIGN**

3

4

## 5    **ABSTRACT**

6    In 2004, Murray et al. published a review of methodological developments in both the design  
7    and analysis of GRTs. In the thirteen years since, there have been many developments in both  
8    areas. The goal of the current paper is to focus on developments in design with a companion  
9    paper to focus on developments in analysis. As a pair, these papers update the 2004 review. This  
10   design paper includes developments in topics included in the earlier review (e.g. clustering,  
11   matching, and individually randomized group treatment trials) and new topics including  
12   constrained randomization and a range of randomized designs that are alternatives to the  
13   standard parallel-arm GRT. These include the stepped wedge GRT, the pseudo-cluster  
14   randomized trial and the network-randomized GRT, which, like the parallel-arm GRT, require  
15   clustering to be accounted for in both their design and analysis.

## 16   **INTRODUCTION**

17   A group-randomized trial (GRT) is a randomized controlled trial in which the unit of  
18   randomization is a group and outcome measurements are obtained on members of those groups.<sup>1</sup>  
19   Also called a cluster randomized trial or community trial,<sup>2-5</sup> a GRT is the best comparative  
20   design available if the intervention operates at a group level, manipulates the physical or social  
21   environment, cannot be delivered to individual members of the group without substantial risk of  
22   contamination across study arms, or if there are other circumstances which warrant the design  
23   such as a desire for herd immunity or a need to estimate both the direct and indirect intervention  
24   effects in studies of infectious diseases.<sup>1-5</sup>

25   In GRTs, outcomes on members of the same group are likely to be more similar to each other  
26   than to outcomes on members from other groups.<sup>1</sup> Such clustering must be accounted for in the

design of GRTs to avoid under-powering the study and accounted for in the analysis to avoid under-estimated standard errors and inflated type I error for the intervention effect.<sup>1-5</sup>

In 2004, Murray et al.<sup>6</sup> published a review of methodological developments in both the design and analysis of GRTs. In the 13 years since, there have been many developments in both areas. The goal of the current paper is to focus on developments in design with a companion paper to focus on developments in analysis.<sup>7</sup> As a pair, these papers update the 2004 review. With both papers, we seek to provide a broad and comprehensive review to guide the reader to seek out appropriate materials for their own circumstances.

## **DEVELOPMENTS IN FUNDAMENTALS OF DESIGN**

### **Clustering**

In its most basic form, a GRT has a hierarchical structure with groups nested within study arm and members nested within groups. Additional levels of nesting may arise through repeated measures over time or from more complex group structures (e.g., children nested in classrooms nested in schools). When designing and analyzing a GRT, it is necessary to account for the clustering associated with the nested design.<sup>1-5</sup>

The intraclass correlation coefficient (ICC) is the clustering measure most commonly used in power calculations and reported in published studies.<sup>8</sup> Eldridge et al.<sup>9</sup> provide a comprehensive review of ICC definitions and measures in general clustered data for both continuous and binary outcomes, the most commonly reported outcomes in GRTs.<sup>10,11</sup> Whereas the ICC for continuous outcome measures is well-defined and generally well understood,<sup>1-4</sup> Eldridge et al.<sup>9</sup> highlight some of the challenges for binary outcomes and provide several definitions (see **Table 1** for the form most commonly presented in GRT texts).<sup>2,4,5,9</sup> Others compare methods to estimate the ICC

of a binary outcome.<sup>12-17</sup> The ICC is not easily defined for rates based on person-time data.<sup>2,4</sup> Recent publications have defined ICC for time-to-event data.<sup>18,19</sup>

The coefficient of variation (CV) is a measure of clustering that is defined for general clustered data when the distributional parameter of interest is a mean, proportion, or rate.<sup>3,17</sup> The CV and ICC for continuous and binary outcomes are related by a mathematical relationship as a function of the distributional parameter of interest (i.e. mean or proportion) and, for continuous outcomes, of the within-group variance,  $\sigma_w^2$  (**Table 1**).<sup>2,4</sup> Hayes and Moulton<sup>2</sup> advocate for the CV generally in power calculations; Donner and Klar agree for event data analyzed as rates.<sup>3</sup>

[TABLE 1 ABOUT HERE]

Given the central role of clustering in planning GRTs, imprecision in the estimated level of clustering can lead to an under-powered trial. Multiple authors address imprecision, and all focus on the ICC.<sup>20-26</sup> Simultaneously, there has been an increasing number of publications that report ICCs (for example, Moerbeek and Teerenstra<sup>27</sup> provide a comprehensive list of such papers) to aid the planning of future studies, consistent with the CONSORT statement on GRTs.<sup>28</sup>

### **Cohort vs. Cross-Sectional GRT Designs**

The choice between a cohort and cross-sectional GRT design (or a combination) is driven by the nature of the research question.<sup>1</sup> The cross-sectional design is preferred when the question is about change in a population<sup>1</sup> or when the time to the outcome is so short as to make a cohort study impractical (e.g., studies involving acute conditions).<sup>2</sup> For example, in order to observe enough participants with malaria at 6-monthly follow-up time points and to be able to draw conclusions about population-level behavior related to malaria treatment choices, Laktabai et al.<sup>29</sup> chose a cross-sectional design in which different population samples were obtained at each

follow-up time point. In contrast, when interested in change in specific individuals, or in mediation, the most natural choice is the cohort design in which a cohort of individuals is enrolled and followed up over time.<sup>1</sup> For example, Turner et al.<sup>30</sup> chose such a design to study child outcomes in mothers with prenatal depression. Similarly, the cohort design is usually required to generate event data in individuals.<sup>2</sup> A combination design could be used whereby the cross-sectional design is augmented by subsampling a cohort of individuals who are followed over time, such as in the COMMIT study.<sup>31</sup> A recent review<sup>32</sup> indicated that the cohort design is the most common GRT design (67% of 75 GRTs).

## **DEVELOPMENTS IN THE DESIGN OF PARALLEL-ARM GROUP-RANDOMIZED TRIALS**

### **Baseline Imbalance of Group Sample Size**

Imbalance of group sample size means that group sizes are different across the groups randomized in the study, with implications for statistical efficiency. Donner discussed variation in group size for GRTs for a design stratified by group size.<sup>33</sup> Guittet et al.<sup>34</sup> and Carter<sup>35</sup> studied the impact on power using simulations, which showed the greatest reduction in power with few groups and/or high ICC. Several authors have offered adjustments to the standard sample size formula for a GRT to correct for variability in group size based on the mean and variance of the group size, or the actual size of each group.<sup>36-39</sup> Others have offered adjustments based on relative efficiency.<sup>40-43</sup> Candel et al.<sup>40,41</sup> reported that relative efficiency ranged from 1.0-0.8 across a variety of distributions for group size with lower values for higher ICCs and greater variability in group size; the minimum relative efficiency was usually no worse than 0.9 for continuous outcomes. They recommended dividing the result from standard formulae for

balanced designs by the relative efficiency for the expected group-size distribution, which was a function of the ICC and the mean and variance of the group size.<sup>40</sup> For binary outcomes, they suggested an additional correction factor based on the estimation method planned for the analysis.<sup>41</sup> You et al.<sup>42</sup> defined relative efficiency in terms of non-centrality parameters; their measure of relative efficiency was a function of the ICC, the mean and variance of the group size, and the number of groups per study arm. Candel and Van Breukelen<sup>43</sup> considered variability not only in group size but also between arms in error variance and the number of groups per arm. They recommended increasing the number of groups in each arm by the inverse of the relative efficiency minus one. Their estimate of the relative efficiency was a function of the number of groups per study arm, the ICC in each study arm, the ratio of the variances in the two study arms, and the mean and variance of the group size.

Consistent across these papers was the recommendation that expectations for variation in group sample size be considered during both the planning stages and the analysis stage. Failure in planning can result in an underpowered study<sup>40-43</sup> while failure in analysis can result in type I error rate inflation.<sup>44</sup>

## **Baseline Imbalance of Covariates**

Imbalance of covariates at baseline threatens the internal validity of the trial. Yet GRTs often randomize a limited number of groups that are heterogeneous in baseline covariates and in baseline outcome measurements. As a result, there is a good chance of baseline covariate imbalance.<sup>6,45</sup> Restricted randomization strategies such as stratification, matching or constrained randomization can be implemented in the design phase to address this issue. However, stratification may have limited use in GRTs if there are more than a handful of covariates to

balance, due to the small number of groups in most trials.<sup>46</sup> Pair-matching also comes with several disadvantages<sup>46</sup> as it affects the proper calculation of ICC<sup>47</sup> and complicates the significance testing of individual-level risk factors.<sup>48</sup> More recently, Imai et al. presented a design-based estimator,<sup>49</sup> which led them to advocate for the use of pair-matching based on the unbiasedness and efficiency of their estimator. Several others highlighted features of this work,<sup>50-52</sup> including the authors' power calculation that does not depend on the ICC, thus avoiding the known ICC problem.<sup>53</sup> Despite efficiency gains of pair-matching over stratification, a simulation study conducted by Imbens led him to conclude that stratified randomization would generally be preferred to pair-matching.<sup>54</sup> We note that strata of size four provide virtually all the advantages of pair-matching while avoiding the disadvantages, and may be preferred over pair-matching for that reason.

To overcome challenges when trying to balance on multiple, possibly continuous, covariates, Raab and Butcher<sup>55</sup> proposed constrained randomization. It is based on a balancing criterion calculated by a weighted sum of squared differences between the study arm means on any group-level or individual-level covariate and seeks to offer better internal validity than both pair-matching and stratification. The approach randomly selects one allocation scheme from a subset of schemes that achieve acceptable balance, identified based on having the smallest values of the balancing criterion. Carter and Hood<sup>56</sup> extended this work to randomize multiple blocks of groups and provided an efficient computer program for public use. The "best balance" score was proposed to measure imbalance of group-level factors under constrained randomization.<sup>57</sup> In simulations with 4 to 20 groups, constrained randomization with the "best balance" score was shown to optimally reduce quadratic imbalances compared with simple randomization, matching and minimization.



Li et al.<sup>58</sup> systematically studied the design parameters of constrained randomization for continuous outcomes, including choice of balancing criterion, candidate set size, and number of covariates to balance. With extensive simulations, they demonstrated that constrained randomization with a balanced candidate subset could improve study power while maintaining the nominal type I error rate, both for a model-based analysis and for a permutation test, as long as the analysis adjusted for potential confounding. Moulton<sup>59</sup> proposed to check for overly constrained designs by counting the number of times each pair of groups received the same study arm allocation. He revealed the risk of inflated type I error in overly constrained designs using a simulation example with 10 groups per study arm. Li et al. further noticed the limitation of overly constrained designs in that they may fail to support a permutation test with a fixed size.<sup>58</sup> In practice, if covariate imbalance is present even after using one of the design strategies described, such imbalance can be accounted for using adjusted analysis that is either pre-planned in the protocol or through post-hoc sensitivity analysis.<sup>7</sup> In summary, constrained randomization seeks to provide both internal validity and efficiency.

## **Methods and Software for Power and Sample Size**

If the ICC is positive, not accounting for it in the analysis will inflate the type I error rate, and the power of the trial will be unknown. If the ICC is estimated as negative, as it can be when the true value is close to zero and sampling error leads to a negative estimate or when there is competition within groups,<sup>1-4,9,60</sup> not accounting for it will reduce the type I error rate so that the test is more conservative, and the power of the trial will be lower than planned.<sup>61</sup> Thus, a good estimate of the ICC is essential for sample size calculation for all GRTs.

One of the simplest power analysis methods often offered for a standard parallel-arm GRT with a single follow-up measurement is to compute the power for an individually randomized trial using the standard formula, and to then inflate this by the design effect,<sup>62</sup> given by  $1 + (m - 1)\rho$ . In this formula,  $m$  is the number of subjects per group and  $\rho$  is the ICC. Unfortunately, this approach only addresses the first of the two penalties associated with group-randomization that were identified by Cornfield almost 40 years ago:<sup>63</sup> extra variation and limited degrees of freedom for the test of the intervention effect. In order to accurately estimate sample size and power for a GRT, it is necessary to also account for the limited degrees of freedom that can arise due to having few groups to randomize. This can be achieved by using appropriate methods detailed in one the GRT texts rather than using the naïve approach of simply inflating the individually randomized trial sample size by the design effect.<sup>1-5,61</sup> In general, appropriate methods calculate sample size using a variance estimate inflated based on the expected ICC and use a t-test rather than a z-test to reflect the desired power and type I error rate, with degrees of freedom based on the number of groups to be randomized.

In practice, both cross-sectional and cohort GRTs are commonly powered based on a comparison between study arms at a single point in time. Then, for GRTs with cohort designs, the analysis section of the study protocol may state that power will be gained by accounting for the repeated measures design in the analysis. However, methods exist for directly computing power in the case of repeated measures in the context of both cross-sectional and cohort designs.<sup>1,27</sup> Authors have noted that regression adjustment for covariates often reduces both the ICC and the residual variance, thereby improving power.<sup>1,64</sup> Heo et al.<sup>65</sup> and Murray et al.<sup>66</sup> provide methods that utilize data from across the entire course of the study, rather than just comparing two means at the end of the study. In practice, the user would require estimates of the variance reduction

expected from repeated measures or from regression adjustment for covariates, which could be obtained from prior studies or pilot data.

Methods exist to power GRTs with additional layers of clustering, whether from additional structural hierarchies<sup>1,67-69</sup> or from the repeated measures in the cohort design.<sup>1,27,64,66,70-73</sup> Konstantopoulos describes how to incorporate cost into the power calculation for three-level GRTs.<sup>74</sup> Hemming et al. discuss approaches to take when the number of groups is fixed ahead of time.<sup>75</sup> Two recent papers focus specifically on binary outcome variables.<sup>13,76</sup> Candel et al. examine the effects of varying group sizes in the context of a two-arm GRT.<sup>77</sup> Durán Pacheco et al. focus on power methods for overdispersed counts.<sup>78</sup> Rutterford et al. and Gao et al. summarize a wide array of methods for sample size calculations in GRTs,<sup>79,80</sup> including for GRT designs involving 1-2 measurements per member or per group and for designs involving 3 or more measurements per member or per group. A new textbook on power analysis for studies with multilevel data also provides a thorough treatment.<sup>27</sup> Previous textbooks on the design and analysis of GRTs devoted at least a chapter to methods for power and sample size.<sup>1-5</sup>

[TABLE 2 ABOUT HERE]

## DEVELOPMENTS IN THE DESIGN OF ALTERNATIVES TO THE PARALLEL-ARM GRT

We discuss four alternatives that can be used in place of a traditional parallel-arm GRT (**Figure 1A, Table 3**). All of these four designs involve randomization and some form of clustering that must be appropriately accounted for in both the design and analysis. As such, they share key

features of the standard parallel-arm GRT yet all have distinct and different features that are important to understand. In practice, some of these designs are still poorly understood.

[TABLE 3 ABOUT HERE]

[FIGURE 1 ABOUT HERE]

### **Stepped Wedge GRT**

The stepped wedge GRT (SW-GRT) is a one-directional crossover GRT in which time is divided into intervals and in which all groups eventually receive the intervention (**Figure 1B**).<sup>81</sup> Systematic reviews indicate increasing popularity.<sup>82-84</sup> Both *Trials* (2015) and the *Journal of Clinical Epidemiology* (2013) recently published special issues focused on the design and analysis of SW-GRTs.

The rationale for this alternative is primarily logistical, i.e., it may not be possible to roll out the intervention in all groups simultaneously,<sup>85-88</sup> although a staggered parallel-arm GRT design could alternatively be used in which blocks of groups were randomized to intervention or control instead of all groups eventually receiving the intervention as in the SW-GRT.<sup>89-91</sup> Others propose a SW-GRT for ethical and acceptability reasons because all groups eventually receive the intervention.<sup>82</sup> This second argument has been discounted as the intervention could be delivered to all control groups at the end of a parallel-arm GRT design,<sup>88,92</sup> often earlier than would be the case in a SW-GRT.<sup>93</sup> When SW-GRTs are conducted in low incidence settings, Hayes et al. emphasized that the order and period of intervention allocation is crucial.<sup>94</sup>

As for the parallel-arm GRT, design choices include cross-sectional<sup>82</sup> vs. cohort<sup>95</sup> with most SW-GRT methodological literature focused on cross-sectional designs whereas most published SW-

GRTs are cohort designs.<sup>96</sup> An additional variation is that of complete vs. incomplete SW-GRTs defined according to whether each group is measured at every time point.<sup>90</sup> Regardless of the specifics of the SW-GRT design, it is important to consider the possible confounding and moderating effects of time in the analysis.<sup>85,90,97-99</sup> Failure to account for both, if they exist, will threaten the internal validity of the study.

Cross-sectional SW-GRT sample size formulae are available for complete and incomplete designs.<sup>90,100-103</sup> Hemming et al. provide a unified approach for the design of both parallel-arm and SW-GRTs and allow for multiple layers of clustering.<sup>90</sup> Cohort SW-GRT sample size calculation relies on simulation.<sup>97,104</sup> Recent work on optimal designs shows that, for large studies, the optimal design is a mixture of a stepped wedge trial embedded in a parallel-arm trial.<sup>105,106</sup> Moerbeek & Teerenstra devote a chapter to sample size methods for SW-GRTs.<sup>27</sup>

### **Network-Randomized GRT**

GRTs have historically been used to minimize the contamination between study arms; such contamination is also called interference.<sup>107</sup> This contamination may give rise to a network of connections between individuals both within- and between-study arms. The latter is of particular relevance to GRT design because it leads to reduced power, although sample size methods exist to preserve power and efficiency.<sup>108</sup>

The network-randomized GRT is a novel design that uses network information to address the challenge of potential contamination in GRTs of infectious diseases.<sup>109-111</sup> In such a design, groups are defined as the network contacts of a disease (index) case and those groups are randomized to study arms. Examples include the snowball trial and the ring trial, each with a distinct way in which the intervention is delivered. In the snowball trial, only the index case

directly receives the intervention, which he is encouraged to share with his contacts (e.g. see Latkin et al.<sup>109</sup> for such a trial of HIV prevention in injection drug users). In the ring trial, ‘rings’ of contacts of the index case are randomized to receive the intervention (**Figure 1C**). This design has been used to study foot-and-mouth,<sup>112</sup> smallpox,<sup>113</sup> and Ebola.<sup>114</sup> For the same sample size, ring trials are more powerful than classical GRTs when the incidence of the infection is low.<sup>115</sup>

### **Pseudo-Cluster Randomized Trial**

In GRTs where all members of the selected groups are recruited to the study, study participants are expected to be representative of the underlying population and, as a result, selection bias is expected to be minimal. In contrast, GRTs with unblinded recruitment after randomization are at risk of selection bias. For example, consider a GRT used to evaluate the effect of a behavioral intervention delivered by providers in the primary care setting. If a provider is first randomized to study arm and then prospectively recruits participants, he may differentially select participants depending on whether he is randomized to the intervention or control arm.<sup>116</sup>

To reduce the risk of such selection bias, Borm et al. introduced the pseudo-cluster randomized trial (PCRT) to allocate intervention to participants in a two-stage process.<sup>117</sup> In the first stage, providers are randomized to a patient allocation-mix (e.g., patients predominantly randomized to intervention vs. patients predominantly randomized to control). In the second stage, patients recruited to the PCRT are individually randomized to intervention or control according to the allocation probability of their provider (e.g., 80% to intervention vs. 20% to intervention) (**Figure 1D**).

An obvious threat to a PCRT design is that the same providers are asked to implement both the intervention and the control arms, depending on which patient they are seeing. Concerns about

contamination are a common reason to randomize providers (i.e. group randomization) so that they deliver either the intervention or the control but not both. The PCRT design would not be appropriate if there are concerns about contamination, and if they exceed concerns about selection bias.

In two published cases, providers were blinded to the two-stage form of randomization and instead assumed that patients were individually randomized to the intervention arm with equal probability.<sup>118,119</sup> Later publications indicate that the PCRT design did well at balancing contamination and selection bias in both studies.<sup>120-122</sup>

Borm et al. provide sample size calculations for continuous outcomes.<sup>117</sup> The clustering by provider (or unit of first stage randomization) must be accounted for in both the design and analysis. No explicit sample size methods are known to be available for non-continuous outcomes. Moerbeek & Teerenstra devote a chapter to sample size methods for PCRTs.<sup>27</sup>

### **Individually Randomized Group Treatment Trial**

Pals et al.<sup>123</sup> identified studies that randomize individuals to study arms but deliver interventions in small groups or through a common change agent as individually randomized group-treatment (IRGT) trials, also called partially clustered or partially nested designs (**Figure 1E**).<sup>72,124</sup> Examples include studies of psychotherapy,<sup>125</sup> weight loss,<sup>126</sup> reduction in sun exposure,<sup>127</sup> and many other outcomes. Clustering associated with these small groups or change agents must be accounted for in the analysis to avoid type I error rate inflation.<sup>72,123,124,128,129</sup> Even so, this accounting appears to be rare in practice.<sup>123,130-133</sup>

Recent papers have reported sample size formulae for IRGT trials with clustering in only one study arm, both for balanced<sup>72,123,128,134</sup> and unbalanced designs.<sup>77,128</sup> Moerbeek & Teerenstra

devote a chapter to sample size methods for IRGT trials focused on methods with clustering in either one or both arms.<sup>27</sup> Roberts addresses sample size methods for IRGT trials in which members belong to more than one small group at the same time or change small groups over the course of the study.<sup>135</sup> Both features have been shown to increase the type I error rate if ignored in the analysis.<sup>135,136</sup>

## **DISCUSSION**

We have summarized many of the most important advances in the design of GRTs during the 13 years since the publication of the earlier review by Murray et al.<sup>6</sup> Many of these developments have focused on alternatives to the standard parallel-arm GRT design, as well as those related to the nature of clustering and its features in all of the designs presented. Space limitations have prevented us from including recent developments involving pilot and feasibility GRTs and group designs such as cutoff designs and regression discontinuity applied to groups. Interested readers are directed to the recently launched *Pilot and Feasibility Studies* peer-reviewed journal and related references<sup>4,137</sup> and to cutoff design references by Pennell et al.<sup>138</sup> and by Schochet.<sup>139</sup>

Through this review, we have sought to ensure that the reader is reminded of the value of good design and gains knowledge in the fundamental principles of a range of recent and potentially beneficial design strategies. Pairing this knowledge with our companion review of developments in the analysis of GRTs,<sup>7</sup> we hope that our work leads to continued improvements in the design and analysis of GRTs.

## **APPENDIX: GLOSSARY**



310 **Balanced candidate subset:** In constrained randomization, where a subset of randomization  
311 schemes is chosen that has “sufficient balance across potentially confounding covariates”  
312 according to “some pre-specified balance metric.”<sup>58</sup>

313 **Baseline covariate balance:** The group-level and individual-level covariate distributions are  
314 similar in all study arms.<sup>55</sup>

315 **Candidate set size:** “The number of possible randomization schemes in a specific  
316 implementation.”<sup>58</sup> “Simple randomization draws from the complete set of candidate schemes,  
317 while constrained randomization considers a subset of schemes.”<sup>58</sup>

318 **Choice of balancing criterion:** Li et al. describe several balancing criteria to assess how well a  
319 GRT is balanced across covariates. These include the “best balance” (BB) metric of de Hoop et  
320 al.,<sup>57</sup> the balance criterion (B) of Raab and Butcher,<sup>55</sup> and the total balance score introduced by  
321 Li et al.<sup>58</sup>

322 **Coefficient of variation:** A measure of between-group variation, defined in Table 1.

323 **Cohort GRT design:** A cohort of individuals is enrolled at baseline and those same individuals  
324 are followed up over time.

325 **Constrained randomization:** Refers “to those designs that go beyond the basic design  
326 constraints to specify classes of randomization outcomes that satisfy certain balancing criteria,  
327 while retaining validity of the design.”<sup>59</sup>

328 **Cross-sectional GRT design:** A different set of individuals is obtained at each time point.

329 **Designed balance at the group level:** When there are equal numbers of groups randomized to  
330 each study arm.

331 **Equivalence:** Assessing whether the new intervention is equivalent to the comparison  
332 intervention.

333 **Individually Randomized Group Treatment Trials:** Studies that randomize individuals to  
334 study arms but deliver treatments in small groups or through a common change agent.<sup>123</sup>

335 **Intraclass correlation:** A measure of between-group variation, defined in Table 1.

336 **Minimization in GRTs:** When the researchers allocate groups to intervention arms based on  
337 groups-specific characteristics in order to achieve a high degree of balance by minimizing the  
338 differences between intervention arms.<sup>57</sup> May be performed sequentially or all at once when  
339 group characteristics are known at the beginning of the study.

340 **Network-Randomized GRT:** The network-randomized GRT is a novel design that uses  
341 network information to address the challenge of potential contamination in GRTs of infectious  
342 diseases.<sup>109-111</sup>

343 **Non-inferiority:** When a trial is designed to show that the new intervention is not worse than  
344 the comparison intervention.

345 **Pair-matching:** At randomization, when groups are matched based on factors thought to be  
346 related to the outcome. Then within each pair of groups, one is allocated at random to one study  
347 arm and the other to the comparison study arm.<sup>140</sup>

348 **Pseudo-cluster randomized trial:** Intervention is allocated to individuals in a two-stage  
349 process. In the first stage, providers are randomized to a patient allocation-mix. In the second  
350 stage, patients recruited to the PCRT are individually randomized to intervention or control  
351 according to the allocation probability of their provider.

**Selection bias:** In some GRTs, groups are randomized before participant recruitment. This can lead to selection bias if researchers (either consciously or unconsciously) recruit specific participants for inclusion in treatment and exclude others based on certain participant characteristics, even when the aforementioned participants are all eligible for participation in the trial (see Farrin et al.<sup>116</sup>).

**Stepped Wedge GRT:** A one-directional crossover GRT in which time is divided into intervals and in which all groups eventually receive the intervention (**Figure 1B**).<sup>81</sup>

**Stratification:** At randomization, when groups are placed into strata based on factors thought to be related to the outcome.<sup>141</sup> Then groups are separately randomized within each strata.

**Superiority:** When a trial is designed to establish whether a new intervention is superior to the comparison intervention (e.g., another drug, a placebo, enhanced usual care). However, the statistical test is still two-sided, allowing for the possibility that the new intervention is actually worse than the comparison.

## **ACKNOWLEDGEMENTS**

To be added later to avoid unblinding during the review process.

## **CONTRIBUTORS**

To be added later to avoid unblinding during the review process.

## **HUMAN PARTICIPANT PROTECTION**

No human subjects participated in this research therefore no IRB approval was sought.

## **Figure 1 The Parallel-Arm GRT and Alternative Group Designs**

Abbreviation: GRT – Group-randomized trial.

Each pictorial representation is an example of the specific design in which baseline measurements are taken. Other versions of each design exist. All examples show 5 individuals per group.

\*The stepped wedge group-randomized trial is a one-directional crossover GRT in which time is divided into intervals and in which all groups eventually receive the intervention, indicated by the shading of the boxes in the figure. The design shown in this figure is known as a “complete design”—that is, every group is measured at every time point. Like parallel-arm GRTs, SW-GRTs can either be cross-sectional or cohort.

†In the PCRT, a group randomized to “intervention” contains a larger proportion of group members receiving the intervention than a group randomized to control.

## REFERENCES

1. Murray DM. *Design and Analysis of Group-Randomized Trials*. New York, NY: Oxford University Press; 1998.
2. Hayes RJ, Moulton LH. *Cluster Randomised Trials*. Boca Raton: CRC Press; 2009.
3. Donner A, Klar N. *Design and Analysis of Cluster Randomization Trials in Health Research*. London: Arnold; 2000.
4. Eldridge S, Kerry S. *A Practical Guide to Cluster Randomised Trials in Health Services Research*. Vol 120: John Wiley & Sons; 2012.
5. Campbell MJ, Walters SJ. *How to Design, Analyse and Report Cluster Randomised Trials in Medicine and Health Related Research*. Chichester, West Sussex: John Wiley & Sons; 2014.
6. Murray DM, Varnell SP, Blitstein JL. Design and analysis of group-randomized trials: a review of recent methodological developments. *Am J Public Health*. 2004;94(3):423-432.
7. Turner EL, Prague M, Gallis JA, Li F, Murray DM. Review of Recent Methodological Developments in Group-Randomized Trials: Part 2 - Analysis. *Am J Public Health*. Submitted.
8. Campbell MK, Fayers PM, Grimshaw JM. Determinants of the intracluster correlation coefficient in cluster randomized trials: the case of implementation research. *Clin Trials*. 2005;2(2):99-107.
9. Eldridge SM, Ukoumunne OC, Carlin JB. The intra-cluster correlation coefficient in cluster randomized trials: a review of definitions. *Int Stat Rev*. 2009;77(3):378-394.
10. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster randomized trials: a systematic review. *Trials*. 2016;17(1):72.
11. Rutterford C, Taljaard M, Dixon S, Copas A, Eldridge S. Reporting and methodological quality of sample size calculations in cluster randomized trials could be improved: a review. *J Clin Epidemiol*. 2015;68(6):716-723.
12. Ridout MS, Demetrio CG, Firth D. Estimating intraclass correlation for binary data. *Biometrics*. 1999;55(1):137-148.
13. Chakraborty H, Moore J, Hartwell TD. Intracluster correlation adjustments to maintain power in cluster trials for binary outcomes. *Contemp Clin Trials*. 2009;30(5):473-480.
14. Thomson A, Hayes R, Cousens S. Measures of between-cluster variability in cluster randomized trials with binary outcomes. *Stat Med*. 2009;28(12):1739-1751.
15. Yelland LN, Salter AB, Ryan P. Performance of the modified Poisson regression approach for estimating relative risks from clustered prospective data. *Am J Epidemiol*. 2011;174(8):984-992.
16. Crespi CM, Wong WK, Wu S. A new dependence parameter approach to improve the design of cluster randomized trials with binary outcomes. *Clin Trials*. 2011;8(6):687-698.
17. Wu S, Crespi CM, Wong WK. Comparison of methods for estimating the intraclass correlation coefficient for binary responses in cancer prevention cluster randomized trials. *Contemp Clin Trials*. 2012;33(5):869-880.
18. Jahn-Eimermacher A, Ingel K, Schneider A. Sample size in cluster-randomized trials with time to event as the primary endpoint. *Stat Med*. 2013;32(5):739-751.
19. Oliveira IR, Molenberghs G, Demétrio CG, Dias CT, Giolo SR, Andrade MC. Quantifying intraclass correlations for count and time-to-event data. *Biom J*. 2016;58(4):852-867.
20. Ukoumunne OC, Davison AC, Gulliford MC, Chinn S. Non-parametric bootstrap confidence intervals for the intraclass correlation coefficient. *Stat Med*. 2003;22(24):3805-3821.
21. Zou G, Donner A. Confidence interval estimation of the intraclass correlation coefficient for binary outcome data. *Biometrics*. 2004;60(3):807-811.
22. Turner RM, Toby Prevost A, Thompson SG. Allowing for imprecision of the intracluster correlation coefficient in the design of cluster randomized trials. *Stat Med*. 2004;23(8):1195-1214.

- 432 23. Turner RM, Thompson SG, Spiegelhalter DJ. Prior distributions for the intracluster correlation  
433 coefficient, based on multiple previous estimates, and their application in cluster randomized  
434 trials. *Clin Trials*. 2005;2(2):108-118.
- 435 24. Turner RM, Omar RZ, Thompson SG. Constructing intervals for the intracluster correlation  
436 coefficient using Bayesian modelling, and application in cluster randomized trials. *Stat Med*.  
437 2006;25(9):1443-1456.
- 438 25. Braschel MC, Svec I, Darlington GA, Donner A. A comparison of confidence interval methods for  
439 the intraclass correlation coefficient in community-based cluster randomization trials with a  
440 binary outcome. *Clin Trials*. 2016;13(2):180-187.
- 441 26. Shoukri MM, Donner A, El-Dali A. Covariate-adjusted confidence interval for the intraclass  
442 correlation coefficient. *Contemp Clin Trials*. 2013;36(1):244-253.
- 443 27. Moerbeek M, Teerenstra S. *Power Analysis of Trials with Multilevel Data*. Boca Raton: CRC Press;  
444 2016.
- 445 28. Campbell MK, Elbourne DR, Altman DG. CONSORT statement: extension to cluster randomised  
446 trials. *Br Med J*. 2004;328(7441):702-708.
- 447 29. Laktabai J, Lesser A, Platt A, et al. An innovative public-private partnership to target subsidized  
448 antimalarials: a study protocol for a cluster randomized controlled trial to evaluate a community  
449 intervention in Western Kenya. In press at BMJ Open.
- 450 30. Turner EL, Sikander S, Bangash O, et al. The effectiveness of the peer delivered Thinking Healthy  
451 Plus (THPP+) Programme for maternal depression and child socio-emotional development in  
452 Pakistan: study protocol for a three-year cluster randomized controlled trial. *Trials*.  
453 2016;17(1):442.
- 454 31. Community Intervention Trial for Smoking Cessation (COMMIT): summary of design and  
455 intervention. *J Natl Cancer Inst*. 1991;83(22):1620-1628.
- 456 32. Murray DM, Pals SP, Blitstein JL, Alfano CM, Lehman J. Design and analysis of group-randomized  
457 trials in cancer: a review of current practices. *J Natl Cancer Inst*. 2008;100(7):483-491.
- 458 33. Donner A. Sample size requirements for stratified cluster randomization designs. *Stat Med*.  
459 1992;11(6):743-750.
- 460 34. Guittet L, Ravaud P, Giraudeau B. Planning a cluster randomized trial with unequal cluster sizes:  
461 practical issues involving continuous outcomes. *BMC Med Res Methodol*. 2006;6:17.
- 462 35. Carter B. Cluster size variability and imbalance in cluster randomized controlled trials. *Stat Med*.  
463 2010;29(29):2984-2993.
- 464 36. Lake S, Kaumann E, Klar N, Betensky R. Sample size re-estimation in cluster randomization trials.  
465 *Stat Med*. 2002;21(10):1337-1350.
- 466 37. Manatunga AK, Hudgens MG, Chen SD. Sample size estimation in cluster randomized studies  
467 with varying cluster size. *Biom J*. 2001;43(1):75-86.
- 468 38. Kerry SM, Bland JM. Unequal cluster sizes for trials in English and Welsh general practice:  
469 implications for sample size calculations. *Stat Med*. 2001;20(3):377-390.
- 470 39. Eldridge SM, Ashby D, Kerry S. Sample size for cluster randomized trials: effect of coefficient of  
471 variation of cluster size and analysis method. *Int J Epidemiol*. 2006;35(5):1292-1300.
- 472 40. van Breukelen G, Candel M, Berger M. Relative efficiency of unequal versus equal cluster sizes in  
473 cluster randomized and multicentre trials. *Stat Med*. 2007;26(13):2589-2603.
- 474 41. Candel MJ, Van Breukelen GJ. Sample size adjustments for varying cluster sizes in cluster  
475 randomized trials with binary outcomes analyzed with second-order PQL mixed logistic  
476 regression. *Stat Med*. 2010;29(14):1488-1501.
- 477 42. You Z, Williams OD, Aban I, Kabagambe EK, Tiwari HK, Cutter G. Relative efficiency and sample  
478 size for cluster randomized trials with variable cluster sizes. *Clin Trials*. 2011;8(1):27-36.

479 43. Candel MJ, Van Breukelen GJ. Repairing the efficiency loss due to varying cluster sizes in two-  
480 level two-armed randomized trials with heterogeneous clustering. *Stat Med*. 2016;35(12):2000-  
481 2015.

482 44. Johnson JL, Kreidler SM, Catellier DJ, Murray DM, Muller KE, Glueck DH. Recommendations for  
483 choosing an analysis method that controls Type I error for unbalanced cluster sample designs  
484 with Gaussian outcomes. *Stat Med*. 2015;34(27):3531-3545.

485 45. Wright N, Ivers N, Eldridge S, Taljaard M, Bremner S. A review of the use of covariates in cluster  
486 randomized trials uncovers marked discrepancies between guidance and practice. *J Clin*  
487 *Epidemiol*. 2015;68(6):603-609.

488 46. Ivers NM, Halperin IJ, Barnsley J, et al. Allocation techniques for balance at baseline in cluster  
489 randomized trials: a methodological review. *Trials*. 2012;13:120.

490 47. Donner A, Klar N. Pitfalls of and controversies in cluster randomized trials. *Am J Public Health*.  
491 2004;26(1):2-19.

492 48. Donner A, Taljaard M, Klar N. The merits of breaking the matches: a cautionary tale. *Stat Med*.  
493 2007;26(9):2036-2051.

494 49. Imai K, King G, Nall C. The essential role of pair matching in cluster-randomized experiments,  
495 with application to the Mexican universal health insurance evaluation. *Stat Sci*. 2009;24(1):29-  
496 53.

497 50. Hill J, Scott M. Comment: The Essential Role of Pair Matching. *Stat Sci*. 2009;24(1):54-58.

498 51. Zhang K, Small DS. Comment: The Essential Role of Pair Matching in Cluster-Randomized  
499 Experiments, with Application to the Mexican Universal Health Insurance Evaluation. *Stat Sci*.  
500 2009;24(1):59-64.

501 52. Imai K, King G, Nall C. Rejoinder: Matched Pairs and the Future of Cluster-Randomized  
502 Experiments. *Stat Sci*. 2009;24(1):65-72.

503 53. Klar N, Donner A. The merits of matching in community intervention trials: a cautionary tale.  
504 *Stat Med*. 1997;16(15):1753-1764.

505 54. Imbens GW. Experimental design for unit and cluster randomized trials. Paper presented at:  
506 Initiative for Impact Evaluation 2011; Cuernavaca, Mexico.

507 55. Raab GM, Butcher I. Balance in cluster randomized trials. *Stat Med*. 2001;20(3):351-365.

508 56. Carter BR, Hood K. Balance algorithm for cluster randomized trials. *BMC Med Res Methodol*.  
509 2008;8:65.

510 57. de Hoop E, Teerenstra S, van Gaal BG, Moerbeek M, Borm GF. The "best balance" allocation led  
511 to optimal balance in cluster-controlled trials. *J Clin Epidemiol*. 2012;65(2):132-137.

512 58. Li F, Lokhnygina Y, Murray DM, Heagerty PJ, DeLong ER. An evaluation of constrained  
513 randomization for the design and analysis of group-randomized trials. *Stat Med*.  
514 2015;35(10):1565-1579.

515 59. Moulton LH. Covariate-based constrained randomization of group-randomized trials. *Clin Trials*.  
516 2004;1(3):297-305.

517 60. Snedecor GW, Cochran WG. *Statistical methods*. 8th ed. Ames: Iowa State University Press;  
518 1989.

519 61. Murray DM, Hannan PJ, Baker WL. A Monte Carlo Study of Alternative Responses To Intraclass  
520 Correlation in Community Trials Is It Ever Possible to Avoid Cornfield's Penalties? *Eval Rev*.  
521 1996;20(3):313-337.

522 62. Donner A, Birkett N, Buck C. Randomization by cluster sample size requirements and analysis.  
523 *Am J Epidemiol*. 1981;114(6):906-914.

524 63. Cornfield J. Randomization by group: a formal analysis. *Am J Epidemiol*. 1978;108(2):100-102.

525 64. Teerenstra S, Eldridge S, Graff M, Hoop E, Borm GF. A simple sample size formula for analysis of  
526 covariance in cluster randomized trials. *Stat Med*. 2012;31(20):2169-2178.

- 527 65. Heo M, Kim Y, Xue X, Kim MY. Sample size requirement to detect an intervention effect at the  
528 end of follow-up in a longitudinal cluster randomized trial. *Stat Med*. 2010;29(3):382-390.
- 529 66. Murray DM, Blitstein JL, Hannan PJ, Baker WL, Lytle LA. Sizing a trial to alter the trajectory of  
530 health behaviours: methods, parameter estimates, and their application. *Stat Med*.  
531 2007;26(11):2297-2316.
- 532 67. Teerenstra S, Lu B, Preisser JS, van Achterberg T, Borm GF. Sample size considerations for GEE  
533 analyses of three-level cluster randomized trials. *Biometrics*. 2010;66(4):1230-1237.
- 534 68. Heo M, Leon AC. Statistical power and sample size requirements for three level hierarchical  
535 cluster randomized trials. *Biometrics*. 2008;64(4):1256-1262.
- 536 69. Teerenstra S, Moerbeek M, van Achterberg T, Pelzer BJ, Borm GF. Sample size calculations for 3-  
537 level cluster randomized trials. *Clin Trials*. 2008;5(5):486-495.
- 538 70. Heo M. Impact of subject attrition on sample size determinations for longitudinal cluster  
539 randomized clinical trials. *J Biopharm Stat*. 2014;24(3):507-522.
- 540 71. Heo M, Leon AC. Sample size requirements to detect an intervention by time interaction in  
541 longitudinal cluster randomized clinical trials. *Stat Med*. 2009;28(6):1017-1027.
- 542 72. Heo M, Litwin AH, Blackstock O, Kim N, Arnsten JH. Sample size determinations for group-based  
543 randomized clinical trials with different levels of data hierarchy between experimental and  
544 control arms. *Stat Methods Med Res*. 2014;26(1):399-413.
- 545 73. Heo M, Xue X, Kim MY. Sample size requirements to detect an intervention by time interaction  
546 in longitudinal cluster randomized clinical trials with random slopes. *Comput Stat Data Anal*.  
547 2013;60:169-178.
- 548 74. Konstantopoulos S. Incorporating cost in power analysis for three-level cluster-randomized  
549 designs. *Eval Rev*. 2009;33(4):335-357.
- 550 75. Hemming K, Girling AJ, Sitch AJ, Marsh J, Lilford RJ. Sample size calculations for cluster  
551 randomised controlled trials with a fixed number of clusters. *BMC Med Res Methodol*.  
552 2011;11:102.
- 553 76. Ahn C, Hu F, Skinner CS, Ahn D. Effect of imbalance and intracluster correlation coefficient in  
554 cluster randomization trials with binary outcomes when the available number of clusters is fixed  
555 in advance. *Contemp Clin Trials*. 2009;30(4):317-320.
- 556 77. Candel MJ, Van Breukelen GJ. Varying cluster sizes in trials with clusters in one treatment arm:  
557 Sample size adjustments when testing treatment effects with linear mixed models. *Stat Med*.  
558 2009;28(18):2307-2324.
- 559 78. Durán Pacheco G, Hattendorf J, Colford JM, Mäusezahl D, Smith T. Performance of analytical  
560 methods for overdispersed counts in cluster randomized trials: Sample size, degree of clustering  
561 and imbalance. *Stat Med*. 2009;28(24):2989-3011.
- 562 79. Rutterford C, Copas A, Eldridge S. Methods for sample size determination in cluster randomized  
563 trials. *Int J Epidemiol*. 2015;44(3):1051-1067.
- 564 80. Gao F, Earnest A, Matchar DB, Campbell MJ, Machin D. Sample size calculations for the design of  
565 cluster randomized trials: A summary of methodology. *Contemp Clin Trials*. 2015;42:41-50.
- 566 81. Spiegelman D. Evaluating public health interventions: 2. Stepping up to routine public health  
567 evaluation with the stepped wedge design. *Am J Public Health*. 2016;106(3):453-457.
- 568 82. Brown CA, Lilford RJ. The stepped wedge trial design: a systematic review. *BMC Med Res*  
569 *Methodol*. 2006;6(1):1.
- 570 83. Mdege ND, Man M-S, Taylor CA, Torgerson DJ. Systematic review of stepped wedge cluster  
571 randomized trials shows that design is particularly used to evaluate interventions during routine  
572 implementation. *J Clin Epidemiol*. 2011;64(9):936-948.
- 573 84. Beard E, Lewis JJ, Copas A, et al. Stepped wedge randomised controlled trials: systematic review  
574 of studies published between 2010 and 2014. *Trials*. 2015;16(1):1-14.



- 575 85. Hussey MA, Hughes JP. Design and analysis of stepped wedge cluster randomized trials.  
576 *Contemp Clin Trials*. 2007;28(2):182-191.
- 577 86. Hargreaves JR, Copas AJ, Beard E, et al. Five questions to consider before conducting a stepped  
578 wedge trial. *Trials*. 2015;16(1):350.
- 579 87. Moulton LH, Golub JE, Durovni B, et al. Statistical design of THRio: a phased implementation  
580 clinic-randomized study of a tuberculosis preventive therapy intervention. *Clin Trials*.  
581 2007;4(2):190-199.
- 582 88. Prost A, Binik A, Abubakar I, et al. Logistic, ethical, and political dimensions of stepped wedge  
583 trials: critical review and case studies. *Trials*. 2015;16(1):351.
- 584 89. Shah More N, Das S, Bapat U, et al. Community resource centres to improve the health of  
585 women and children in Mumbai slums: study protocol for a cluster randomized controlled trial.  
586 *Trials*. 2013;14:132.
- 587 90. Hemming K, Lilford R, Girling AJ. Stepped-wedge cluster randomised controlled trials: a generic  
588 framework including parallel and multiple-level designs. *Stat Med*. 2015;34(2):181-196.
- 589 91. Kotz D, Spigt M, Arts ICW, Crutzen R, Viechtbauer W. Use of the stepped wedge design cannot  
590 be recommended: A critical appraisal and comparison with the classic cluster randomized  
591 controlled trial design. *J Clin Epidemiol*. 2012;65(12):1249-1252.
- 592 92. Kotz D, Spigt M, Arts IC, Crutzen R, Viechtbauer W. Researchers should convince policy makers  
593 to perform a classic cluster randomized controlled trial instead of a stepped wedge design when  
594 an intervention is rolled out. *J Clin Epidemiol*. 2012;65(12):1255.
- 595 93. Murray DM, Pennell M, Rhoda D, Hade EM, Paskett ED. Designing studies that would address  
596 the multilayered nature of health care. *J Natl Cancer Inst Monogr*. 2010(40):90-96.
- 597 94. Hayes RJ, Alexander ND, Bennett S, Cousens SN. Design and analysis issues in cluster-  
598 randomized trials of interventions against infectious diseases. *Stat Methods Med Res*.  
599 2000;9(2):95-116.
- 600 95. Copas AJ, Lewis JJ, Thompson JA, Davey C, Baio G, Hargreaves JR. Designing a stepped wedge  
601 trial: three main designs, carry-over effects and randomisation approaches. *Trials*.  
602 2015;16(1):352.
- 603 96. Hemming K, Haines TP, Chilton PJ, Girling AJ, Lilford RJ. The stepped wedge cluster randomised  
604 trial: rationale, design, analysis, and reporting. *BMJ*. 2015;350:h391.
- 605 97. Baio G, Copas A, Ambler G, Hargreaves J, Beard E, Omar RZ. Sample size calculation for a  
606 stepped wedge trial. *Trials*. 2015;16(1):354.
- 607 98. Handley MA, Schillinger D, Shiboski S. Quasi-experimental designs in practice-based research  
608 settings: design and implementation considerations. *The Journal of the American Board of*  
609 *Family Medicine*. 2011;24(5):589-596.
- 610 99. Liao X, Zhou X, Spiegelman D. A note on "Design and analysis of stepped wedge cluster  
611 randomized trials". *Contemp Clin Trials*. 2015;45(Pt B):338-339.
- 612 100. Hemming K, Taljaard M. Sample size calculations for stepped wedge and cluster randomised  
613 trials: a unified approach. *J Clin Epidemiol*. 2016;69:137-146.
- 614 101. Hemming K, Girling A. A menu-driven facility for power and detectable-difference calculations in  
615 stepped-wedge cluster-randomized trials. *Stata J*. 2014;14(2):363-380.
- 616 102. Hughes J. Calculation of power for stepped wedge design. <http://tinyurl.com/hwp5dgr>.  
617 Accessed January 12, 2017.
- 618 103. Hughes J. Calculation of power for stepped wedge design (means). <http://tinyurl.com/jvcr5bu>.  
619 Accessed January 12, 2017.
- 620 104. Baio G. SWSamp: Simulation-based sample size calculations for a Stepped Wedge Trial (and  
621 more). 2016; <https://sites.google.com/a/statistica.it/gianluca/swsamp>.

622 105. Lawrie J, Carlin JB, Forbes AB. Optimal stepped wedge designs. *Stat Probab Lett.* 2015;99:210-  
623 214.

624 106. Girling AJ, Hemming K. Statistical efficiency and optimal design for stepped cluster studies under  
625 linear mixed effects models. *Stat Med.* 2016;35(13):2149-2166.

626 107. Hudgens MG, Halloran ME. Toward causal inference with interference. *J Am Stat Assoc.*  
627 2008;103(482):832-842.

628 108. Wang R, Goyal R, Lei Q, Essex M, De Gruttola V. Sample size considerations in the design of  
629 cluster randomized trials of combination HIV prevention. *Clin Trials.* 2014;11(3):309-318.

630 109. Latkin C, Donnell D, Liu TY, Davey-Rothwell M, Celentano D, Metzger D. The dynamic  
631 relationship between social norms and behaviors: the results of an HIV prevention network  
632 intervention for injection drug users. *Addiction.* 2013;108(5):934-943.

633 110. Staples PC, Ogburn EL, Onnela J-P. Incorporating Contact Network Structure in Cluster  
634 Randomized Trials. *Sci Rep.* 2015;5:17581.

635 111. Harling G, Wang R, Onnela J, De Gruttola V. Leveraging contact network structure in the design  
636 of cluster randomized trials. *Clin Trials.* 2016 [Epub ahead of print].

637 112. Keeling M, Woolhouse M, May R, Davies G, Grenfell B. Modelling vaccination strategies against  
638 foot-and-mouth disease. *Nature.* 2003;421(6919):136-142.

639 113. Kretzschmar M, Van den Hof S, Wallinga J, Van Wijngaarden J. Ring vaccination and smallpox  
640 control. *Emerg Infect Dis.* 2004;10(5):832-841.

641 114. Enserink M. High hopes for Guinean vaccine trial. *Science.* 2015;347(6219):219-220.

642 115. Henao-Restrepo AM. The ring vaccination trial: a novel cluster randomised controlled trial  
643 design to evaluate vaccine efficacy and effectiveness during outbreaks, with special reference to  
644 Ebola. *BMJ.* 2015;351.

645 116. Farrin A, Russell I, Torgerson D, Underwood M. Differential recruitment in a cluster randomized  
646 trial in primary care: the experience of the UK back pain, exercise, active management and  
647 manipulation (UK BEAM) feasibility study. *Clin Trials.* 2005;2(2):119-124.

648 117. Borm GF, Melis RJ, Teerenstra S, Peer PG. Pseudo cluster randomization: a treatment allocation  
649 method to minimize contamination and selection bias. *Stat Med.* 2005;24(23):3535-3547.

650 118. Melis RJ, van Eijken MI, Borm GF, et al. The design of the Dutch EASYcare study: a randomised  
651 controlled trial on the effectiveness of a problem-based community intervention model for frail  
652 elderly people. *BMC Health Serv Res.* 2005;5:65.

653 119. Pence BW, Gaynes BN, Thielman NM, et al. Balancing contamination and referral bias in a  
654 randomized clinical trial: An application of pseudo-cluster randomization. *Am J Epidemiol.*  
655 2015;182(12):1039-1046.

656 120. Melis RJ, Teerenstra S, Rikkert MG, Borm GF. Pseudo cluster randomization performed well  
657 when used in practice. *J Clin Epidemiol.* 2008;61(11):1169-1175.

658 121. Pence BW, Gaynes BN, Adams JL, et al. The effect of antidepressant treatment on HIV and  
659 depression outcomes: results from a randomized trial. *AIDS.* 2015;29(15):1975-1986.

660 122. Teerenstra S, Melis R, Peer P, Borm G. Pseudo cluster randomization dealt with selection bias  
661 and contamination in clinical trials. *J Clin Epidemiol.* 2006;59(4):381-386.

662 123. Pals SP, Murray DM, Alfano CM, Shadish WR, Hannan PJ, Baker WL. Individually randomized  
663 group treatment trials: a critical appraisal of frequently used design and analytic approaches.  
664 *Am J Public Health.* 2008;98(8):1418-1424.

665 124. Baldwin SA, Bauer DJ, Stice E, Rohde P. Evaluating models for partially clustered designs. *Psychol*  
666 *Methods.* 2011;16(2):149-165.

667 125. Carlbring P, Bohman S, Brunt S, et al. Remote treatment of panic disorder: a randomized trial of  
668 internet-based cognitive behavior therapy supplemented with telephone calls. *Am J Psychiatry.*  
669 2006;163(12):2119-2125.

126. Jeffery RW, Linde JA, Finch EA, Rothman AJ, King CM. A Satisfaction Enhancement Intervention for Long-Term Weight Loss. *Obesity*. 2006;14(5):863-869.
127. Jackson KM, Aiken LS. Evaluation of a multicomponent appearance-based sun-protective intervention for young women: uncovering the mechanisms of program efficacy. *Health Psychol*. 2006;25(1):34.
128. Roberts C, Roberts SA. Design and analysis of clinical trials with clustering effects due to treatment. *Clin Trials*. 2005;2(2):152-162.
129. Kahan BC, Morris TP. Assessing potential sources of clustering in individually randomised trials. *BMC Med Res Methodol*. 2013;13:58.
130. Pals SL, Wiegand RE, Murray DM. Ignoring the group in group-level HIV/AIDS intervention trials: a review of reported design and analytic methods. *AIDS*. 2011;25(7):989-996.
131. Lee KJ, Thompson SG. Clustering by health professional in individually randomised trials. *BMJ*. 2005;330(7483):142-144.
132. Biau DJ, Porcher R, Boutron I. The account for provider and center effects in multicenter interventional and surgical randomized controlled trials is in need of improvement: a review. *J Clin Epidemiol*. 2008;61(5):435-439.
133. Oltean H, Gagnier JJ. Use of clustering analysis in randomized controlled trials in orthopaedic surgery. *BMC Med Res Methodol*. 2015;15:17.
134. Moerbeek M, Wong WK. Sample size formulae for trials comparing group and individual treatments in a multilevel model. *Stat Med*. 2008;27(15):2850-2864.
135. Roberts C, Walwyn R. Design and analysis of non-pharmacological treatment trials with multiple therapists per patient. *Stat Med*. 2013;32(1):81-98.
136. Andridge RR, Shoben AB, Muller KE, Murray DM. Analytic methods for individually randomized group treatment trials and group-randomized trials when subjects belong to multiple groups. *Stat Med*. 2014;33(13):2178-2190.
137. Eldridge SM, Costelloe CE, Kahan BC, Lancaster GA, Kerry SM. How big should the pilot study for my cluster randomised trial be? *Stat Methods Med Res*. 2016;25(3):1039-1056.
138. Pennell ML, Hade EM, Murray DM, Rhoda DA. Cutoff designs for community-based intervention studies. *Stat Med*. 2011;30(15):1865-1882.
139. Schochet PZ. Statistical power for regression discontinuity designs in education evaluations. *J Educ Behav Stat*. 2009;34(2):238-266.
140. Campbell M, Donner A, Klar N. Developments in cluster randomized trials and Statistics in Medicine. *Stat Med*. 2007;26(1):2-19.
141. Kernan WN, Viscoli CM, Makuch RW, Brass LM, Horwitz RI. Stratified randomization for clinical trials. *J Clin Epidemiol*. 1999;52(1):19-26.

**Table 1. Two Common Measures of Clustering for General Clustered Data for Two Common Types of Outcome**

Outcome measure	Intra-cluster correlation coefficient (ICC, $\rho$ ) <sup>a</sup>	Coefficient of Variation (CV, $k$ )	Relationship of ICC to CV <sup>b</sup>
Continuous	$\sigma_B^2 / (\sigma_B^2 + \sigma_W^2)$	$\sigma_B / \mu$	$1 / \left( 1 + \frac{\sigma_W^2}{k^2 \mu^2} \right)$
Binary	$\sigma_B^2 / \pi(1 - \pi)$	$\sigma_B / \pi$	$k^2 \pi / (1 - \pi)$

Note:  $\mu$  = overall mean for continuous outcome data;  $\pi$  = overall proportion for binary outcome data;  $\sigma_B^2$  = between-group variance;  $\sigma_W^2$  = within-group variance (i.e. residual error variance). As is common practice, the two clustering measures are for general clustered data and do not focus on the GRT design in which the intervention effect is of primary interest (e.g. see Chapter 2 of Hayes and Moulton<sup>2</sup> for more details). The intervention parameter of interest in GRT is typically: difference of means for continuous outcomes; difference of proportions, ratio of proportions or odds ratio for binary outcomes; rate difference or rate ratio for event outcomes.

<sup>a</sup> There are multiple definitions of the ICC for binary outcomes (see <sup>12-17</sup>). The specific formulation provided here is one of the simplest and most commonly used (see, for example, equation (2.4) of Hayes and Moulton<sup>2</sup> and equation (8) of Eldridge et al.<sup>9</sup>).

<sup>b</sup> Note that, while the relationship for binary outcomes is only a function of  $k$  and the distributional parameter of interest ( $\pi$ ), the relationship for continuous outcomes is a function of both the distributional parameter of interest ( $\mu$ ) and  $\sigma_W^2$ .

716 **Table 2. Software for Sample Size Calculations in Parallel-Arm GRTs**

Software	Functionality
PASS	Sample size calculations for GRTs comparing two means (non-inferiority, equivalence, or superiority), two proportions (non-inferiority, equivalence, or superiority), two Poisson rates, and for a logrank test.
nQuery	Comparison of two means, proportions, and rates.
Stata	User-provided command <code>clustersampsi</code> . Can compute sample size for continuous, binary, and rate outcomes for two-sided tests in equal-sized arms.
R	Package <code>CRTSize</code> for comparing two means or two binary proportions.
SAS	No built-in functionality at this time.
Calculator	For some simple designs, parameter values can be plugged in to formulas provided in textbooks and online.

717

718

719 **Table 3. Characteristics of the Parallel-Arm Group Randomized Trial (GRT) and of Alternative Group Designs**

Design	Acronym	One-stage randomization		Two-stage randomization	Type of follow-up possible <sup>20</sup>	
		By Group	By Individual		Cross-sectional	Cohort
<b>Parallel-Arm GRT</b>	GRT	✓	-	-	✓	✓
<b>Stepped Wedge GRT</b>	SW-GRT	✓	-	-	✓	✓
<b>Network-Randomized GRT</b>	NR-GRT	✓	-	-	-	✓ <sup>1</sup>
<b>Pseudo-Cluster Randomized Trial</b>	PCRT	-	-	✓	-	✓ <sup>2</sup>
<b>Individually Randomized Group Treatment Trial</b>	IRGT trial	-	✓	-	-	✓ <sup>3</sup>

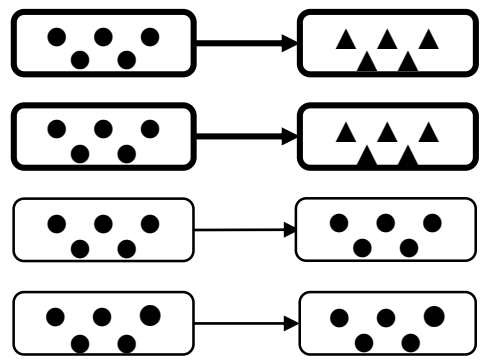
<sup>1</sup> In the network-randomized GRT, the index case and its network is usually defined at baseline and therefore the design is expected to use a cohort design and not allow a cross-sectional design

<sup>2</sup> In the pseudo-cluster randomized trial, because randomization is undertaken in two stages with individuals randomized to intervention or control in the second stage, the design requires that a cohort of individuals be enrolled at study baseline in order to be followed over time

<sup>3</sup> In the individually randomized group treatment trial, individual randomization is performed and therefore, like the pseudo-cluster randomized trial, a cohort of individuals is enrolled and followed over time.

Baseline Follow-up

A: Parallel GRT

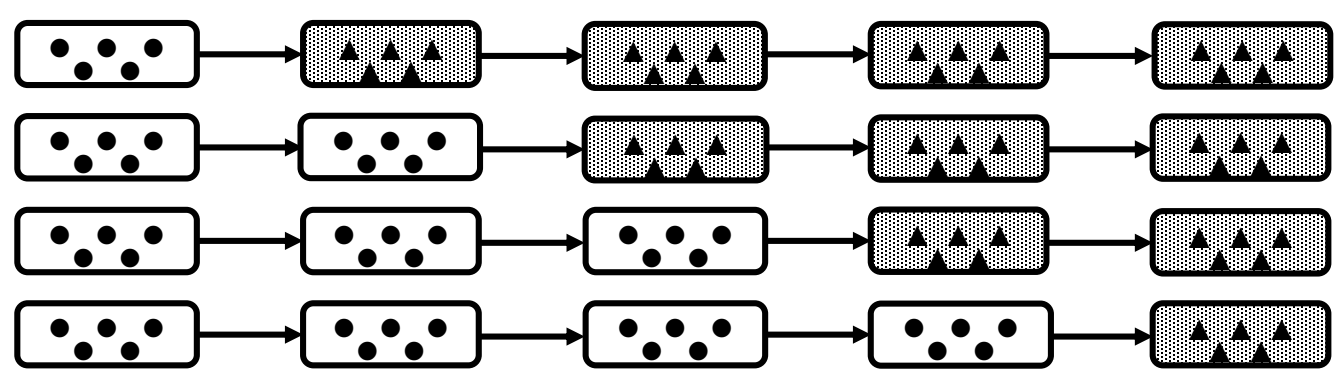


Groups are randomized to intervention or control at baseline, then either the same individuals are followed up over time (cohort GRT) or different individuals in the same group are sampled at different time points (cross-sectional GRT).

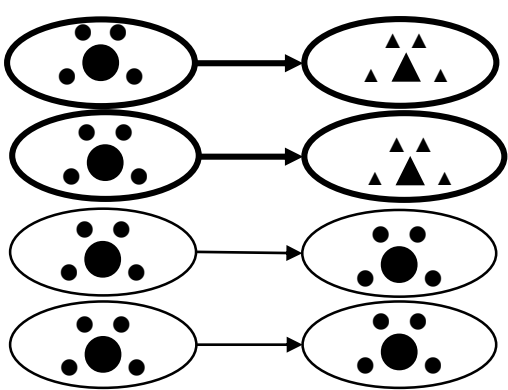
Legend

- ▲ Individual measured under intervention
- Individual measured under no intervention
- Group randomized to intervention
- Group randomized to control

B: Stepped Wedge GRT (SW-GRT)\*

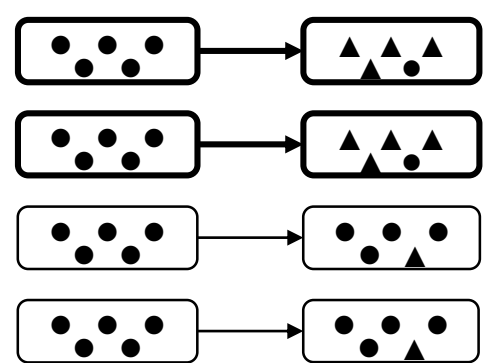


C: Network-Randomized GRT (NR-GRT)



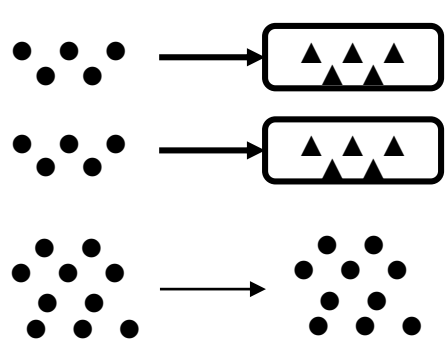
Groups are defined as the network contacts of an index disease case and those groups of contacts are then randomized to intervention or control. The larger symbols represent the index case in each group.

D: Pseudo-Cluster Randomized Trial (PCRT)†



Assignment to intervention is based on a two-stage process. In the first stage, groups (e.g., providers) are randomized to a patient allocation-mix, here shown as predominantly (80%) intervention vs. predominantly (80%) control. In the second stage, patients recruited to the PCRT are individually randomized to intervention or control.

E: Individually Randomized Group Treatment (IGRT) Trial



Individuals are randomized to intervention or control but treatments are delivered in small groups or through a common change agent.